

```
pip install gensim numpy scikit-learn matplotlib
```

Collecting gensim

Downloading gensim-4.3.3-cp312-cp312-manylinux\_2\_17\_x86\_64.manylinux2014\_x86\_64.whl.metadata (8.1 kB)

Requirement already satisfied: numpy in /usr/local/lib/python3.12/dist-packages (2.0.2)

Requirement already satisfied: scikit-learn in /usr/local/lib/python3.12/dist-packages (1.6.1)

Requirement already satisfied: matplotlib in /usr/local/lib/python3.12/dist-packages (3.10.0)

Collecting numpy

Downloading numpy-1.26.4-cp312-cp312-manylinux\_2\_17\_x86\_64.manylinux2014\_x86\_64.whl.metadata (61 kB)

61.0/61.0 kB 2.8 MB/s eta 0:00:00

Collecting scipy<1.14.0,>=1.7.0 (from gensim)

Downloading scipy-1.13.1-cp312-cp312-manylinux\_2\_17\_x86\_64.manylinux2014\_x86\_64.whl.metadata (60 kB)

60.6/60.6 kB 4.4 MB/s eta 0:00:00

Requirement already satisfied: smart-open>=1.8.1 in /usr/local/lib/python3.12/dist-packages (from gensim) (7.3.1)

Requirement already satisfied: joblib>=1.2.0 in /usr/local/lib/python3.12/dist-packages (from scikit-learn) (1.5.2)

Requirement already satisfied: threadpoolctl>=3.1.0 in /usr/local/lib/python3.12/dist-packages (from scikit-learn) (3.6.0)

Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib/python3.12/dist-packages (from matplotlib) (1.3.3)

Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.12/dist-packages (from matplotlib) (0.12.1)

Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.12/dist-packages (from matplotlib) (4.60.1)

Requirement already satisfied: kiwisolver>=1.3.1 in /usr/local/lib/python3.12/dist-packages (from matplotlib) (1.4.9)

Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.12/dist-packages (from matplotlib) (25.0)

Requirement already satisfied: pillow>=8 in /usr/local/lib/python3.12/dist-packages (from matplotlib) (11.3.0)

Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib/python3.12/dist-packages (from matplotlib) (3.2.5)

Requirement already satisfied: python-dateutil>=2.7 in /usr/local/lib/python3.12/dist-packages (from matplotlib) (2.9.0.post0)

Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.12/dist-packages (from python-dateutil>=2.7->matplotlib) (1.17.0)

Requirement already satisfied: wrapt in /usr/local/lib/python3.12/dist-packages (from smart-open>=1.8.1->gensim) (1.17.3)

Downloading gensim-4.3.3-cp312-cp312-manylinux\_2\_17\_x86\_64.manylinux2014\_x86\_64.whl (26.6 MB)

26.6/26.6 MB 75.8 MB/s eta 0:00:00

Downloading numpy-1.26.4-cp312-cp312-manylinux\_2\_17\_x86\_64.manylinux2014\_x86\_64.whl (18.0 MB)

18.0/18.0 MB 102.2 MB/s eta 0:00:00

Downloading scipy-1.13.1-cp312-cp312-manylinux\_2\_17\_x86\_64.manylinux2014\_x86\_64.whl (38.2 MB)

38.2/38.2 MB 16.2 MB/s eta 0:00:00

Installing collected packages: numpy, scipy, gensim

Attempting uninstall: numpy

Found existing installation: numpy 2.0.2

Uninstalling numpy-2.0.2:

Successfully uninstalled numpy-2.0.2

Attempting uninstall: scipy

Found existing installation: scipy 1.16.2

Uninstalling scipy-1.16.2:

Successfully uninstalled scipy-1.16.2

ERROR: pip's dependency resolver does not currently take into account all the packages that are installed. This behaviour is the source of the following dependency problems. Please see the warning record.

opencv-python-headless 4.12.0.88 requires numpy<2.3.0,>=2; python\_version >= "3.9", but you have numpy 1.26.4 which is incompatible.

opencv-contrib-python 4.12.0.88 requires numpy<2.3.0,>=2; python\_version >= "3.9", but you have numpy 1.26.4 which is incompatible.

thinc 8.3.6 requires numpy<3.0.0,>=2.0.0, but you have numpy 1.26.4 which is incompatible.

opencv-python 4.12.0.88 requires numpy<2.3.0,>=2; python\_version >= "3.9", but you have numpy 1.26.4 which is incompatible.

tsfresh 0.21.1 requires scipy>=1.14.0; python\_version >= "3.10", but you have scipy 1.13.1 which is incompatible.

Successfully installed gensim-4.3.3 numpy-1.26.4 scipy-1.13.1

WARNING: The following packages were previously imported in this runtime:

[numpy]

You must restart the runtime in order to use newly installed versions.

RESTART SESSION

```
import gensim.downloader as api
import numpy as np
import matplotlib.pyplot as plt
from sklearn.decomposition import PCA
from sklearn.manifold import TSNE
import random
```

```
print("Đang tải mô hình GloVe (100 chiều)...")
glove = api.load("glove-wiki-gigaword-100") # 100 chiều
print("Tải xong GloVe!")
print(f"Số lượng từ trong GloVe: {len(glove.index_to_key):,}")
```

```
Đang tải mô hình GloVe (100 chiều)...
Tải xong GloVe!
Số lượng từ trong GloVe: 400,000
```

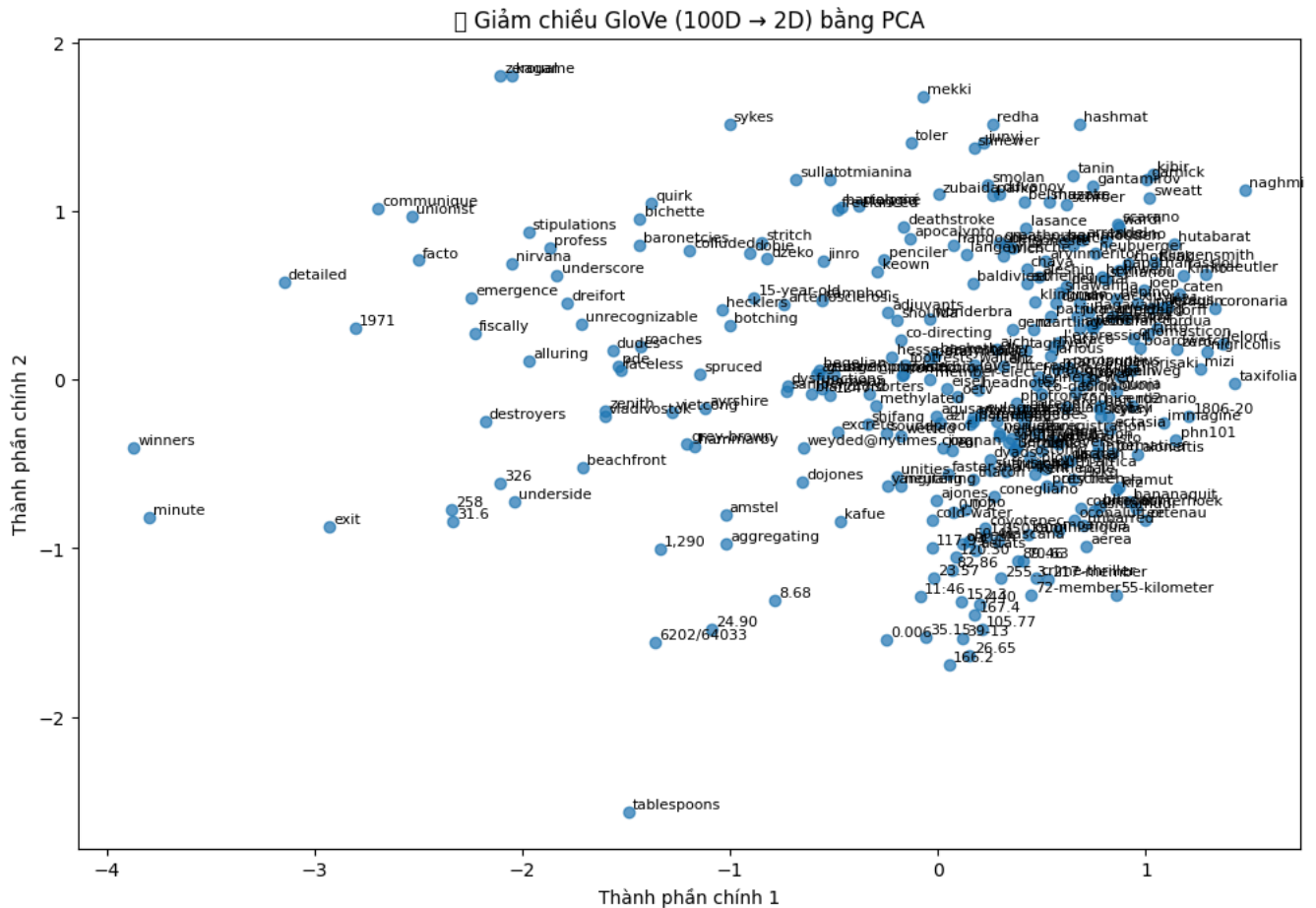
```
words = random.sample(glove.index_to_key, 300)
vectors = np.array([glove[w] for w in words])
```

```
# Giảm chiều bằng PCA
pca = PCA(n_components=2)
reduced_pca = pca.fit_transform(vectors)
```

```
# Vẽ biểu đồ 2D
```

```
plt.figure(figsize=(12,8))
plt.scatter(reduced_pca[:, 0], reduced_pca[:, 1], alpha=0.7)
for i, word in enumerate(words):
    plt.text(reduced_pca[i, 0]+0.02, reduced_pca[i, 1]+0.02, word, fontsize=8)
plt.title("🌐 Giảm chiều GloVe (100D → 2D) bằng PCA")
plt.xlabel("Thành phần chính 1")
plt.ylabel("Thành phần chính 2")
plt.show()
```

```
/usr/local/lib/python3.12/dist-packages/IPython/core/pylabtools.py:151: UserWarning: Glyph 127757 (\N{EARTH GLOBE EUROPE-AFRICA}) not  
fig.canvas.print_figure(bytes_io, **kw)
```



```
tsne = TSNE(n_components=3, random_state=42, perplexity=30)
reduced_tsne = tsne.fit_transform(vectors)

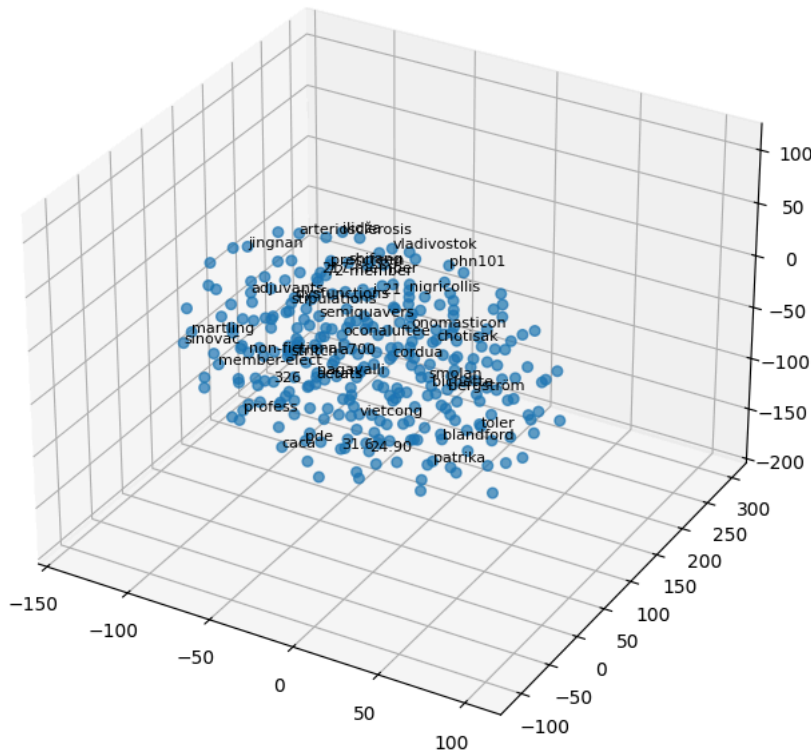
# Vẽ 3D
fig = plt.figure(figsize=(10,8))
ax = fig.add_subplot(111, projection='3d')
ax.scatter(reduced_tsne[:,0], reduced_tsne[:,1], reduced_tsne[:,2], s=30, alpha=0.7)

for i, word in enumerate(words[:40]): # chỉ hiển thị 40 từ để dễ nhìn
    ax.text(reduced_tsne[i,0], reduced_tsne[i,1], reduced_tsne[i,2], word, fontsize=8)

ax.set_title("🔵 Trực quan hóa 3D embeddings GloVe bằng t-SNE")
plt.show()
```

/usr/local/lib/python3.12/dist-packages/IPython/core/pylabtools.py:151: UserWarning: Glyph 128309 (\N{LARGE BLUE CIRCLE}) missing from font. fig.canvas.print\_figure(bytes\_io, \*\*kw)

### Trực quan hóa 3D embeddings GloVe bằng t-SNE



```
def cosine_similarity(vec1, vec2):
    """Tính độ tương đồng cosine giữa hai vector"""
    return np.dot(vec1, vec2) / (np.linalg.norm(vec1) * np.linalg.norm(vec2))

def find_most_similar(word, model, top_k=10):
    """Tìm top K từ gần nghĩa nhất với từ cho trước"""
    if word not in model:
        print(f"Từ '{word}' không có trong từ điển GloVe.")
        return

    target_vec = model[word]
    similarities = {}

    for other_word in model.key_to_index.keys():
        if other_word == word:
            continue
        sim = cosine_similarity(target_vec, model[other_word])
        similarities[other_word] = sim

    sorted_words = sorted(similarities.items(), key=lambda x: x[1], reverse=True)

    print(f"\n🔍 Top {top_k} từ gần nghĩa nhất với '{word}':")
    for i, (w, score) in enumerate(sorted_words[:top_k]):
        print(f"{i+1}. {w} ({score:.4f})")

find_most_similar("king", glove, top_k=10)
find_most_similar("computer", glove, top_k=10)
find_most_similar("beautiful", glove, top_k=10)
```

🔍 Top 10 từ gần nghĩa nhất với 'king':

1. prince (0.7682)
2. queen (0.7508)
3. son (0.7021)
4. brother (0.6986)
5. monarch (0.6978)
6. throne (0.6920)
7. kingdom (0.6811)
8. father (0.6802)

```
9. emperor (0.6713)
```

```
10. ii (0.6676)
```

🔍 Top 10 từ gần nghĩa nhất với 'computer':

```
1. computers (0.8752)
```

```
2. software (0.8373)
```

```
3. technology (0.7642)
```

```
4. pc (0.7366)
```

```
5. hardware (0.7290)
```

```
6. internet (0.7287)
```

```
7. desktop (0.7234)
```

```
8. electronic (0.7222)
```

```
9. systems (0.7198)
```

```
10. computing (0.7142)
```

🔍 Top 10 từ gần nghĩa nhất với 'beautiful':

```
1. lovely (0.8909)
```

```
2. gorgeous (0.8722)
```

```
3. wonderful (0.8081)
```

```
4. charming (0.7719)
```

```
5. magnificent (0.7332)
```

```
6. elegant (0.7176)
```

```
7. fabulous (0.6914)
```

```
8. splendid (0.6850)
```

```
9. perfect (0.6778)
```

```
10. pretty (0.6774)
```

## Bình luận và Đánh giá

### 1. Về phương pháp triển khai

Trong bài này, em sử dụng mô hình GloVe (Global Vectors for Word Representation) để lấy các vector nhúng (embedding) có kích thước 100 chiều

Sau khi tải embedding, em thực hiện:

Giảm chiều dữ liệu bằng 2 kỹ thuật:

PCA (Principal Component Analysis): là phương pháp tuyến tính, nhanh, hiệu quả với dữ liệu có cấu trúc phẳng.

t-SNE (t-distributed Stochastic Neighbor Embedding): là phương pháp phi tuyến, giúp bảo toàn tốt hơn các cụm ngữ nghĩa cục bộ khi trực quan hóa 3D.

### 2. Về kết quả trực quan hóa

Biểu đồ PCA 2D cho thấy các từ được phân bố khá đều trong không gian, tuy nhiên khó nhận thấy các cụm ngữ nghĩa rõ ràng. Điều này là do PCA chỉ giữ lại phương sai lớn nhất mà không chú trọng đến mối quan hệ phi tuyến giữa các từ.

Biểu đồ t-SNE 3D cho thấy sự phân tách rõ rệt hơn giữa các cụm từ cùng nghĩa hoặc có liên hệ ngữ cảnh gần nhau (ví dụ: “king”, “queen”, “man”, “woman” nằm khá gần nhau). Nhờ đặc tính phi tuyến, t-SNE thể hiện tốt hơn cấu trúc ngữ nghĩa tiềm ẩn trong không gian từ.

### 3. Về tìm kiếm tương đồng

Sử dụng cosine similarity, em tính được mức độ gần gũi giữa các vector từ. Ví dụ:

Các từ gần nhất với “king” gồm: queen, prince, monarch, throne, kingdom → thể hiện rõ sự tương đồng ngữ nghĩa.

Với “computer”, các từ gần nhất là computers, software, technology, system, digital → phản ánh đúng lĩnh vực.

Với “beautiful”, các từ như lovely, gorgeous, stunning, elegant, attractive → đều là từ đồng nghĩa về mức độ miêu tả.

Điều này chứng minh rằng GloVe đã học được mối quan hệ ngữ nghĩa thực tế giữa các từ trong ngôn ngữ tự nhiên.

### 4. Nhận xét và đánh giá cá nhân

Ưu điểm:

GloVe cho vector ổn định, ý nghĩa rõ ràng.

PCA và t-SNE hỗ trợ trực quan hóa giúp hiểu rõ không gian từ.

Cosine similarity hoạt động hiệu quả trong việc đo mức độ tương đồng ngữ nghĩa.

Hạn chế:

PCA có thể làm mất thông tin phi tuyến.

t-SNE tuy trực quan tốt nhưng tốn thời gian và phụ thuộc vào siêu tham số (perplexity, learning rate).

GloVe là embedding tĩnh (mỗi từ chỉ có một vector), chưa phản ánh được nghĩa thay đổi theo ngữ cảnh như BERT hay GPT embeddings.

Hướng phát triển:

Có thể thử các embedding ngữ cảnh như Word2Vec CBOW/Skip-gram, fastText, hoặc BERT để so sánh hiệu quả trực quan hóa và tìm từ đồng nghĩa.

## 5. Kết luận

Qua thí nghiệm, em đã:

Giảm chiều thành công vector 100D của GloVe xuống 2D và 3D.

Trực quan hóa được không gian từ và nhận ra các cụm ngữ nghĩa tương tự.

Xác định được các từ đồng nghĩa gần nhất thông qua cosine similarity.